# COMMSCOPE®

# Cabling considerations of AI data centers

**Dr. Earl Parsons, Director Data Center/Intelligent Building Architecture Evolution**

## Introduction

For decades the danger of malicious artificial intelligence (AI) has been a trope in science fiction. Film antagonists like HAL 9000, the Terminator, the Replicants, and the robots from the Matrix are opposing forces to the plucky humans who must overcome the dangers of technology. Recently, the release of DALLE-2 and ChatGPT has captured the imagination of the wider public of what AI can do. This has led to discussions on how AI will change the nature of education and work.

AI is the main driver for current and future data center growth. There are three aspects to AI:

- During training, a large set of data is fed into the algorithm, allowing it to learn.
- Inference AI then takes information and analyzes it. For example, is this a picture of a cat?
- Generative AI is the most exciting because, from simple prompts, the algorithm can output text or images that have never been created before.

The computation for AI is carried out in graphical processing units (GPUs). These specialized chips are best at parallel processing and are well suited to AI. These models used to train and run AI are too large for a single machine. Figure 1 shows the historical growth of AI models in PetaFLOPs (floating point operations). Multiple GPUs, spread over many servers and racks, are required to handle these large models. These GPUs need to be connected to allow them to do the work of AI, and this paper outlines the challenges and opportunities of cabling AI data centers—targeted primarily to large enterprises, who will build their own AI clusters—and instructs them on the best way to cable their AI clusters.
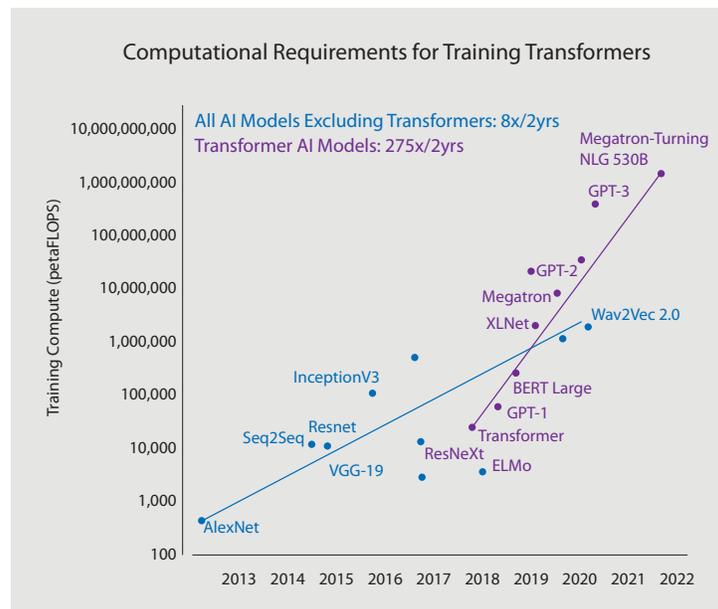


Fig. 1: AI model size in petaFLOPS
(source: https://blogs.nvidia.com/blog/2022/03/25/what-is-a-transformer-model/)

## Typical data hall architecture

Nearly all modern data centers, especially hyperscale data centers, use a Folded Clos architecture, also called leaf-spine. All the leaf switches in a data center connect to all the spine switches. In the data hall, server racks connect to a top-of-rack (ToR) switch. The ToR is then connected to a leaf switch at the end of the row or in another room with fiber cable. The servers in the rack are connected to the ToR with short copper cables. These copper cables are one to two meters long and carry 25G or 50G signaling.

This configuration uses few fiber cables in the data hall. For example, Meta data centers that use the F16 architecture (see Figure 2) will have 16 duplex fiber cables from each of the server racks in a row. These cables run from the ToR to the end of the row, where they connect with modules that combine duplex fibers to 24 fiber cables. The 24 fiber cables then run to another room to connect to leaf switches.
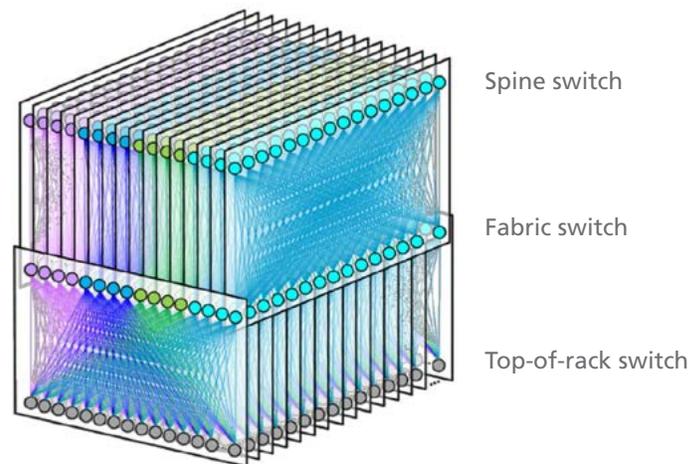


Fig. 2: FaceBook F16 data center network topology
(source: https://engineering.fb.com/2019/03/14/data-center-engineering/f16-minipack/)

Data centers that implement AI will house AI clusters next to compute clusters with traditional architecture. Traditional compute is sometimes called the front-end network, and the AI clusters are sometimes called the back-end network.

## Data halls with AI clusters

AI clusters require a new data center architecture. The GPU servers require much more connectivity between servers, but there are fewer servers per rack due to power and heat restraints. This leads to situations where we have more inter-rack cabling than traditional data centers. Each GPU server is connected to a switch within the row or room. These links require 100G to 400G at distances that cannot be supported by copper. In addition, each server requires connectivity to the switch fabric, storage, and out-of-band management.

## Example: NVIDIA

As an example, we can look at the architecture proposed by NVIDIA, a leader in the AI space. NVIDIA's latest GPU server is the DGX H100 and has 4x800G ports to switches (operated as 8x400GE), 4x400GE ports to storage, and 1GE and 10GE ports for management. A DGX SuperPOD (as in Figure 3) can contain 32 of these GPU servers connected to 18 switches in a single row. Each row would then have 384x400GE fiber links for switch fabric and storage and 64 copper links for management. This is a remarkable increase in the number of fiber links in the data hall. The F16 architecture mentioned above would have 128 (8x16) duplex fiber cables with the same number of server racks.

Fig. 3: Rendering of NVIDIA SuperPOD
(source: https://www.nvidia.com/en-us/data-center/dgx-superpod/)

## What link lengths are in an AI cluster?

In the ideal scenario illustrated by NVIDIA, all the GPU servers in an AI cluster will be close together. AI/machine learning algorithms, like high-performance computing, are extremely sensitive to link latency. One estimate claimed that 30 percent of the time to run a large training model was spent on network latency and 70 percent was spent on compute time. Since training a large model can cost up to $10 million, this networking time represents a significant cost. Even a latency saving of 50 nanoseconds, or 10 m of fiber, is significant. Nearly all the links in AI clusters are limited to 100 m reaches.

Unfortunately, not all data centers will be able to locate the GPU server racks in the same row. These racks require ~40 kilowatts to power the GPU servers. This is more power than typical server racks, and data centers built with lower power requirements will need to space out their GPU racks.

## Which transceivers should you use?

Operators should carefully consider which optical transceivers and fiber cables they will use in their AI clusters to minimize cost and power consumption. As explained above, the longest links within an AI cluster will be limited to 100 m. Due to the short reach, the optics cost will be dominated by the transceiver. Transceivers that use parallel fiber will have an advantage: they do not require the optical multiplexers and demultiplexers used for wavelength division multiplexing. This results in both lower cost and lower power for transceivers with parallel fiber. The transceiver cost savings more than offset the small increase in cost for a multifiber cable instead of a duplex fiber cable. For example, 400G-DR4 transceivers with eight fiber cables is more cost effective than 400G-FR4 transceivers with duplex fiber cable.

Links up to 100 m are supported by singlemode fiber and multimode fiber applications. Advances like silicon photonics have reduced the cost of singlemode transceivers to bring them closer to the cost of equivalent multimode transceivers. Our market research indicates that, for high-speed transceivers (400G+), the cost of a singlemode transceiver is double the cost of an equivalent multimode transceiver. While multimode fiber has a slightly higher cost than singlemode fiber, the difference in cable cost between multimode and singlemode is smaller since multifiber cable costs are dominated by MPO connectors.

In addition, high-speed multimode transceivers use one to two watts less power than their singlemode counterparts. With 768 transceivers in a single AI cluster (128 memory links + 256 switch links X2) a setup using multimode fiber will save up to 1.5 kW. This may seem small compared to the 10 kW that each DGX H100 consumes, but for AI clusters any opportunity to save power will be welcome.

In IEEE 802.3db a new multimode transceiver was standardized: the VR or very short reach. This application targets in-row cabling like AI clusters with max reach of 50 m. These transceivers have the potential to offer the lowest cost and power consumption for AI connectivity.

## Transceivers vs. AOCs

Many AI/ML clusters and HPCs use active optical cables (AOCs) to interconnect GPUs and switches. An active optical cable is a fiber cable with integrated optical transmitters and receivers on either end. Most AOCs are used for short reaches and typically use multimode fiber and VCSELs. High-speed (>40G) active optical cables will use the same OM3 or OM4 fiber used in fiber cables that connect optical transceivers. The transmitters and receivers in an AOC may be the same as in analogous transceivers but are the castoffs; each transmitter and receiver doesn't need to meet rigorous interoperability specs; they only need to operate with the specific unit attached to the other end of the cable. Since no optical connectors are accessible to the installer, the skills required to clean and inspect fiber connectors are not needed.

The downside of AOCs is that they do not have the flexibility offered by transceivers. Installing AOCs is time-consuming as the cable must be routed with the transceiver attached. Correctly installing AOCs with breakouts is especially challenging. The failure rate for AOCs is double that of equivalent transceivers. When an AOC fails, a new AOC must be routed through the network. This takes away from the compute time. Finally, when it is time to upgrade the network links, the AOCs must be removed and replaced with new AOCs. With transceivers, the fiber cabling is part of the infrastructure and may remain in place for several generations of data rates.

## Conclusions

Careful consideration of the cabling of AI clusters will help save cost, power, and installation time. The right fiber cabling will enable organizations to fully benefit from artificial intelligence.

CommScope pushes the boundaries of communications technology with game-changing ideas and ground-breaking discoveries that spark profound human achievement. We collaborate with our customers and partners to design, create and build the world's most advanced networks. It is our passion and commitment to identify the next opportunity and realize a better tomorrow. Discover more at commscope.com.

**COMMSCOPE®**